

Editorial Manager(tm) for PLoS ONE
Manuscript Draft

Manuscript Number:

Title: Identification of Antifreeze Proteins and Their Functional Residues by Support Vector Machine and Genetic Algorithms based on n-Peptide Compositions

Short Title: Identify AFPs and Their Functional Residues

Article Type: Research Article

Section/Category: Other

Keywords: support vector machines; genetic algorithm; n-peptide composition; antifreeze protein; AFP

Corresponding Author: Chin Sheng Yu, Ph.D.

Corresponding Author's Institution: Feng Chia University

First Author: Chin Sheng Yu, Ph.D.

Order of Authors: Chin Sheng Yu, Ph.D.;Chih Hao Lu

Abstract: For the first time, multiple sets of global n-peptide compositions from antifreeze protein (AFP) sequences of certain cold-adapted fish and insects were analyzed using support vector machine and genetic algorithms. The identification of AFPs is difficult because they exist as evolutionarily divergent types, and because their sequences and structures are present in limited numbers in currently available databases. Our results reveal that it is feasible to identify the shared sequential features among the various structural types of AFPs. Moreover, we were able to identify residues involved in ice binding without referring to three-dimensional structures of AFPs. This approach should be useful for genomic and proteomic studies involving cold-adapted organisms.

Suggested Reviewers: Peter L. Davies
Queen's University
peter.davies@queensu.ca
expert of antifreeze protein

Brendan J J McConkey
University of Waterloo
mcconkey@uwaterloo.ca
expert of antifreeze protein recognition

Opposed Reviewers:

Dear Prof.,

Here within enclosed is our paper for consideration to be published on PloS ONE.

The further information about the paper is in the following:

The Title: **Identification of antifreeze proteins and their important residues
by using support vector machines based on *n*-peptide
compositions**

The Authors: Chin-Sheng Yu and Chih-Hao Lu

It is first discussed that the antifreeze proteins and their functional important residues can be identified from protein sequences analysis. The common characters in antifreeze sequence still lack due to the poor homologs and radical different type in current database. Our approach not only provides excellent results for discriminating them without using the 3D structural information, but the most important, it is allowed a further investigation the rule of potential key residues in ice-binding interface.

The authors claim that none of the material in the paper has been published or is under consideration for publication elsewhere.

I am the corresponding author and my address and other information is as follows,

Address: Department of Information Engineering and Computer Science,
Feng Chia University, Taichung, 40724, Taiwan

E-mail: yucs@fcu.edu.tw

Tel: 886-4-24517250 ext. 3742

Fax: 886-4-24516101

Thank you very much for consideration!

Identification of Antifreeze Proteins and Their Functional Residues by Support Vector Machine and Genetic Algorithms based on n -Peptide Compositions

Chin-Sheng Yu^{1,2*} and Chih-Hao Lu³

From the ¹Department of Information Engineering and Computer Science, ²Master's Program in Biomedical Informatics and Biomedical Engineering, Feng Chia University, Taichung 40724, Taiwan and the ³Graduate Institute of Molecular Systems Biomedicine, China Medical University, Taichung 40402, Taiwan

*Correspond to: Chin-Sheng Yu, Department of Information Engineering and Computer Science, Feng Chia University, Taichung 40724, Taiwan. FAX: +886-4-2451-6101. Phone: +886-4-2451-7250, ext. 3742. E-mail: yucs@fcu.edu.tw.

1 Abstract

2 For the first time, multiple sets of global n -peptide compositions from antifreeze protein (AFP)
3 sequences of certain cold-adapted fish and insects were analyzed using support vector machine and
4 genetic algorithms. The identification of AFPs is difficult because they exist as evolutionarily
5 divergent types, and because their sequences and structures are present in limited numbers in currently
6 available databases. Our results reveal that it is feasible to identify the shared sequential features
7 among the various structural types of AFPs. Moreover, we were able to identify residues involved in
8 ice binding without referring to three-dimensional structures of AFPs. This approach should be useful
9 for genomic and proteomic studies involving cold-adapted organisms.

10 Keywords: support vector machines; genetic algorithm; n -peptide composition; antifreeze protein;
11 AFP

12 INTRODUCTION

13 Antifreeze proteins (AFPs) in cold-adapted organisms prevent macroscopic ice build-up by binding to
14 ice and thereby forestalling additional crystallization [1]. By doing so, AFPs allow organisms to
15 survive below 0°C. It is of great interest to harness this singular property—non-antifreeze proteins
16 cannot bind ice—for applications related to the agriculture and food industries [2,3,4,5] and to the
17 rational design of new AFPs. However, first it is necessary to understand how AFPs and ice interact.
18 Accurately identifying AFPs from evolutionarily divergent organisms is difficult because their
19 sequences and structures differ radically [6,7]. To complicate matters further, for closely related
20 species, the sequences, and consequently the structures, of their AFPs may also differ substantially if
21 they have been geographically isolated [8]. Additionally, searching for homologous sequences within
22 databases has not been a fruitful approach given the disparity among AFP sequences. Directly

23 studying AFP-ice interactions is also difficult, and a definitive picture of such interactions is not
24 currently available [7]. Therefore, because many AFPs do not have structural or sequential features in
25 common, it is challenging to correlate the relationships among their sequences, structures, and
26 function.

27 A large number of biochemical and structural studies [9,10,11] have been performed in an attempt to
28 understand how AFPs interact with ice on the molecular level, including site-directed mutagenesis
29 [12,13,14] and computational experiments [15]. An ice-binding model that incorporates surface
30 complementarity is generally accepted [16]. Recently, Doxey and colleagues [9] successfully
31 identified AFPs, for which three-dimensional (3D) crystallographic structures were available, on the
32 basis of their highly ordered and planar ice-binding surfaces, but their algorithm could not identify an
33 AFP when only its NMR solution structure was available because the coordinates for the atoms at and
34 near its surface were not well defined. [9,17]. Additionally, their algorithm requires the use of a
35 three-dimensional (3D) structure, which is not always available for a given AFP.

36 It is obvious, therefore, that AFPs cannot be easily distinguished from other types of proteins.
37 Additional information is needed to understand how AFPs and ice interact on a fundamental
38 physicochemical level before such interactions can be applied to cold-adapted mechanisms. Although
39 the types of amino acids present are closely coupled to the ice-binding properties of AFPs [10,13],
40 current models usually rely on only 3D structures. To make additional use of the knowledge that has
41 accumulated over the decades, e.g., identification of the "hydrophobic surface" effect [7,11], the
42 spatial regularity of an AFP solvent accessible surface, the presence of nonpolar residues, and other
43 properties directly related to the binding properties of AFPs, an algorithm that can discern these
44 properties is necessary. Therefore, for this report, we developed an integrated approach to rapidly
45 identify AFPs from their amino acid sequences. Our statistically based, support vector machine (SVM)
46 algorithm has been used to identify certain inherent protein traits e.g., protein disulfide connectivities
47 [18], subcellular localizations [19,20], and protein folds [21], when given a query sequence, and it
48 does not require a computational mechanical model or structure comparison. For this report, during
49 the training and testing of this algorithm for different classifiers associated with AFPs, multiple
50 feature schemes based on n -peptide compositions extracted from the sequences were used. Then, a
51 genetic algorithm (GA) was used iteratively for key-feature selection and to improve the identification
52 accuracy. This integrated approach enabled the recognition of AFPs on the basis of preferred short
53 peptide sequences, rather than on structural comparisons. The identified AFP sequence features have
54 not been reported previously, yet they correlate well with the properties of the ice-binding interfaces.
55 This approach is suitable for the further identification of the ice-binding surfaces of AFPs.

56 **METHODS**

57 *The Validation Dataset that Contained AFPs and non-AFPs with Known 3D Structures—*

58 To assess our approach without bias, we tested it using a sequence validation dataset that did not

59 contain homologous proteins, and to examine the effects of key residues on function, we included
 60 only AFPs that had solved structures. This set contained 3762 nonredundant non-AFPs and 44 AFPs,
 61 which had been collected from the PISCES server [22] and the Protein Data Bank (PDB) [23],
 62 respectively. To include as many representative structures as possible, the non-AFPs had <25%
 63 pairwise sequence identity (SI), R-factors of 0.25 and a crystallographic resolution of at least 2 Å. The
 64 AFP sequences were separated into eight subsets on the basis of sequence identity by ClustalW2 [24].
 65 Table 1 lists the PDB IDs of the AFPs in each subset. For a given subset, the associated AFP(s) had a
 66 sequence(s) that was not homologous to any of the AFPs in the other subsets. The non-AFPs were
 67 randomly divided among the eight subsets to cross test the performance of our approach and then
 68 were merged as a single trained model for use with other (independent) datasets (see below). Under
 69 such a critical condition, any afterward AFPs recognition so far is not simply from the self-trained
 70 sequences.

71

72 *Independent Datasets—*

73 We constructed three other datasets that did not contain the AFPs included in the aforementioned eight
 74 subsets to test our algorithm after training it with the latter. The first set included three AFP structures
 75 deposited recently in the PDB [23]; the second set contained 369 nonredundant AFP sequences
 76 deposited in the UniProKT database [25,26], which represented an evolutionarily divergent group of
 77 organisms; the third set contained two “antifreeze-like” (AFL) proteins that, while incapable of
 78 binding ice, have both a sequence and a structure that are very similar to the fish type III AFP [27].
 79 Table 2 lists the number of AFPs derived from each type of organism included in the second dataset.

80

81 *Feature schemes—*

82 The n -peptide composition feature-based coding schemes, with $n = 1$ encoding the amino acid
 83 composition; $n = 2$, the dipeptide composition; $n = 3$, the tripeptide composition, etc., were used
 84 previously to predict protein properties [19,20,21,28], and we used them to characterize the important
 85 ice-binding features of AFPs. A set of symbols, A_n for the original amino acids; H_n for hydrophobicity
 86 [29]; V_n for the normalized van der Waals volume [29]; Z_n for polarizability [29]; P_n for polarity [29];
 87 and F_n , S_n , and E_n , for groups of residues classified according to four, seven, and eight
 88 physical/chemical properties, respectively, were used to denote the feature schemes [19]. However, to
 89 characterize the key functional residues more robustly, partitioned subsequences, g -gap dipeptides,
 90 and local amino acid composition strategies were also included. [19] The partitioned amino acid
 91 composition X_k^Y is a concatenation of all amino acid sequences of composition Y and length k . The
 92 symbol D_g identifies the frequency of a sequence in the form $a(x)_g b$, where a and b denote specific
 93 amino acids and $(x)_g$ denotes the g -intervening (g -gap) residues of any type between the pair. The
 94 symbol W_l indicates the amino acid composition for peptides characterized by a set of sliding
 95 windows of length l centered on a given type of amino acid. It provides information concerning the
 96 sequential neighbors for of a given type of amino acid.

97

98 *Assembly Machine-learning Algorithms—*

99 All SVM calculations were performed using LIBSVM [30], which is a general library for support
 100 vector classification and regression, and the radial basis function kernel. In addition to the SVM
 101 algorithm [31], we implemented a GA to efficiently optimize the selection of feature attributes as
 102 detailed previously [18]. The combined use of the SVM algorithm and the GA is denoted as SVMGA.
 103 For the SVMGA, the feature attributes of each feature scheme, the penalty parameter C, the kernel
 104 parameter γ of the RBF function used for SVM identification by the GA approach were determined in
 105 advance. The GA procedure rapidly filtered out feature attributes that are not useful for SVM
 106 identification on the basis of each feature scheme.

107

108 *The Voting System—*

109 The coding scheme symbols given above denote the SVM classifiers that were derived from the
 110 various properties of the sequence descriptors. For simplicity, the participants in the
 111 SVM-identification system [19,20] were incorporated as:

$$112 \sum_{k=1}^9 X_k^{A_i} + \sum_{g=0}^6 D_g + \sum_S X_{k=5}^S + \sum_{l \in S'} W_l$$

113 with $S = \{H_3, V_3, Z_3, P_3, F_3, S_2, E_2\}$ and $S' = \{7, \dots, 15\}$. The system counts the jury votes from each
 114 classifier to determine if a protein is an AFP.

115

116 *Performance Assessment—*

117 As in previous work [19,20,21], we employed the accuracy $Q_i = c_i/n_i \times 100$ to assess the performance
 118 of identification, i.e., the prediction accuracy, where c_i is the number of correctly identified AFPs in
 119 the class $i \in (\text{AFP}, \text{non-AFP})$, and n_i is the number of sequences. The overall identification accuracy is
 120 given by

$$121 P = \sum_i f_i Q_i,$$

122 where $f_i = n_i/N$, and N is the total number of sequences. Although Q_i provides a convenient assessment
 123 for identification performance, the Matthews Correlation Coefficient (*MCC*) [32] is a more
 124 informative measure of the performance and is given by:

$$125 MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FN)(TP + FP)(TN + FP)(TN + FN)}},$$

126 where TP , TN , FP , and FN are the number of true positives, true negatives, false positives, and false
 127 negatives, respectively. A value for *MCC* of 1, 0, or -1 represents a perfect correlation, a random
 128 correlation, or an inverse correlation, respectively. Consideration of the *MCC*, allowed us to modify
 129 our approach to lower the number of false positives returned. To be a credible method, our approach
 130 needed to return as few false positives as possible.

131

132 *AFP Sequence Homology Search—*

133 To verify our ability to identify AFPs via their protein sequences, we tested the homology
134 relationships among the AFP sequences. A query sequence from the second independent data set was
135 aligned with the sequences of the 44 AFPs of the validation set. Only these 44 AFPs were used
136 because their 3D structures have been solved, and they had been experimentally shown to bind ice.
137 We performed an all-against-all sequence alignment using the global alignment program ALIGN [33].
138 Only the top-ranked sequence of the 44 AFP sequences was then used to assess the effect of homology
139 on AFP identification, i.e., the SI value for the query sequence and the top-ranked sequence
140 determined the usefulness of the homology search approach.

141 RESULTS**142** *Identification of AFPs in a Cross-validation Dataset—*

143 For the cross-validation test, the non-AFPs were randomly and equally divided into eight subsets,
144 each of which contained a single representative AFP (which is identified by the first PDB ID (in bold
145 type) in each subset list (Table 1)), and these sets formed the single representative AFP mode. Then, if
146 the AFP representative had homologous sequences, these sequences were added into the
147 corresponding subset. The eight subsets can be thought of as eight distant branches of an evolutionary
148 tree. These sets formed the multiple representative AFP mode. For an experiment, the sequences of
149 seven of the subsets were used to train the SVM algorithm with a given feature scheme, and then the
150 output model of the trained algorithm was used to test the sequences in the subset that was not used
151 for training. This training-and-testing cross-validation procedure was repeated eight times for a given
152 feature scheme, each time omitting a different sequence subset during training. All results reported the
153 performance on the total number of datasets. The SVM classifiers were optimized so that the
154 algorithm could assign a protein sequence as either an AFP or non-AFP sequence.

155 Table 3 contains a summary of the identification accuracies and the MCC values for the different
156 combinations of feature schemes used for the single representative AFP mode and the multiple
157 representative AFP mode. Only the best result for a given feature scheme is reported. The best overall
158 identification accuracy was 62.5% for the single representative AFP mode used by the SVM
159 algorithm. Incorporation of the GA algorithm substantially improved the identification accuracy.
160 Using the iterative procedures mentioned above, the GA identified the largest number of true positives
161 and the smallest number of false positives as it discarded feature attributes that were not useful for the
162 SVM classification. The assembled SVMGA approach correctly identified all AFPs in the single
163 representative AFP mode. Using just the smallest possible number of selected features, the SVM
164 classifier identified more completely structurally dissimilar AFPs than did Doxey and colleagues who
165 used the structural characteristics of the AFPs [9]. After we decreased the number of FPs as much as
166 possible (<70 FPs remained), we tested the algorithm with the multiple representative AFP mode,
167 which was a more realistic dataset. Although the performance of the algorithm declined with the

168 increase in the number of divergent sequences, the identification accuracy was a respectable 54.5%.

169

170 *Identification of AFPs in the Independent Datasets—*

171 The three AFPs of the first independent dataset, which were the A chains of 2zib, 3bog, and 3boi were
 172 all accurately identified as AFPs. We observed that the sequence of 2zib is homologous to that of 2afp,
 173 which was contained in the eighth validation subset, and the sequences of 3bog and 3boi are
 174 homologous to that of 2pne, which was contained in the sixth validation subset. In addition to
 175 accurately identifying the proteins of the first independent dataset as AFPs, the algorithm also
 176 recognized that the human and bacterial AFL proteins (PDB IDs 1wvo and 1xuz, respectively) [27]
 177 were not AFPs. The human AFL and the bacterial AFL are both very similar in sequence and structure
 178 to that of the fish type III AFP (PDB code 1msi).

179 For the AFPs of the second independent dataset, which represent a divergent group of organisms and
 180 were collected from the UniProKT database [25,26], about 61% were correctly identified as AFPs by
 181 the SVMGA. The SI pair distribution, which characterizes the relative number of sequence pairs in
 182 the close percentage sequence identity interval, was used to examine the effect of sequence homology
 183 on AFP identification. The 369 AFP sequences were each used as a query sequence to profile the SI
 184 pair-distribution. Each query sequence was aligned with the 44 AFPs of the validation set and also
 185 with the other 368 sequences of the second independent data set. The largest SI value for each query
 186 that was aligned with the 44 AFPs was plotted along the y axis, and the largest SI value for
 187 corresponding sequence aligned with the other 368 sequences of the second dataset was plotted along
 188 the x axis (Fig. 1). The SI values associated with AFPs in the independent dataset that were
 189 incorrectly identified by the SVMGA are colored red in Figure 1, and most of these values are <20%,
 190 which below the so-called midnight-zone threshold where a structural/functional relationship can be
 191 detected [34]. Because the dataset that contained the 369 AFPs was biased as it contained AFPs from
 192 well-characterized cold-adapted organisms, many of the points were located at the far end of the x
 193 axis.

194

195 *Coding Schemes—*

196 For the different coding-scheme SVM classifiers used in this study, we were able to reduce the
 197 number of feature attributes required by at least 50% after implementing the GA. Consequently, each
 198 remaining classifier was well suited to identifying the corresponding type of AFP (Table 4). To
 199 understand why the features were selected as classifiers, we assigned a number (vote) when the
 200 pattern of residues in a sequence matched a GA-selected feature attribute of a coding scheme. The
 201 sequence position was marked as an SVMGA key residue if it had received a majority of the jury
 202 votes from the 14 coding schemes that we used for the multiple representative AFP mode. For
 203 instance, the dipeptide LT was selected in the D_0 scheme, and the interval dipeptide T(X_2)T was
 204 selected in the D_2 scheme. Hence, for the short peptide NTALT, the L in the fourth position and the
 205 first T each received one vote, and the second T received two votes (Table 5). Eight representative

206 AFPs are presented in Fig. 2, with their SVMGA key residues marked. Residues with >6 votes, with 4
207 or 5 votes, and with <3 votes are colored red, yellow, and gray, respectively. Fig. 3 illustrates the
208 average number of SVMGA key residues in AFP sequences (black bars) and the number of in
209 non-AFP sequences (gray bars). And it is obviously that the number of SVMGA key residues in AFP
210 sequences is twice in non-AFPs. Approximately 70% of the SVMGA-selected key residues are
211 solvent exposed (data not shown), which is sensible as these residues are more likely to interact with
212 ice.

213 DISCUSSION

214 Previous studies have deduced the structural character of the interactions between ice and AFP
215 molecules [7,14]. Knowing how ice and AFP molecules interact allows for the identification of AFPs
216 given their structures (see the excellent results of Doxey and colleagues reported in Table 3). However,
217 the method of Doxey and colleagues required the use of proteins with solved 3D structures, and
218 therefore, until this report, there has not been a more general method for AFP identification.

219 For this report, we presented an integrated machine-learning method, SVMGA, to identify AFPs that
220 uses multiple *n*-peptide composition features. Our results show that sequentially divergent AFPs can
221 be identified according to their shared sequence characteristics because any test sequence or its
222 homologs are not appearing in trained set. A set of *n*-peptide composition-based SVM predictors were
223 combined to accurately recognize AFPs, and more importantly, to identify the key functional residues
224 at the ice-binding surfaces. Several reports [7] have characterized defining residue repeats in AFP
225 sequences, e.g., alanine-rich sequences in the α -helix of type I AFPs (A28–A34, Fig. 2f), and
226 Thr-Cys-Thr (Fig. 2b) or Thr-Xaa-Thr (Fig. 2c) sequences in insect AFPs. The feature attributes,
227 selected by our SVMGA approach, included these features. Some of the key SVMGA residues in
228 these representative structures of AFPs, formed relatively flat planes, e.g., the red and yellow
229 clustered regions in Fig. 2 and 4. Additionally, SVMGA approach identified some residues reside at
230 the interface between two chains of crystallized form in PDB, e.g., T13 and T24 in chain A of 1wfa
231 (Fig. 2f), but actually the active protein is monomer. We found others that the SVMGA key residues
232 in red, L12, L23, A31, and T35, reside on the same side of the flat binding interface. Another similar
233 example is the β -sheet plane of chain A in 1ezg (Fig. 2b), although the Thr-Cys-Xaa tri-peptide
234 parallel strands [35] align perfectly in the dimer crystallized form, this flattest ice-binding surface is
235 found in the monomer as seen by the coloration at the functional interface.

236 We also inspected the key residues that were identified in the eelpout type III AFP, which has been
237 subjected to many mutagenesis studies. As mentioned in Method, this eelpout type III AFP, which
238 PDB codes 1msi, had no homolog in any of the AFPs in trained subsets 1, 3, 4, 5, 6, 7 and 8 (Table 1.).
239 And the key residues of 1msi were inferred from these dissimilar trained sequences by SVMGA
240 approach. Compared with previous studies [12,14], the SVMGA identified half of the proven
241 ice-binding residues at the interface (Fig. 4b). For the three residues, N14, A16, and T18, which when

242 mutated caused the greatest decreases in AFP activity, the SVMGA method found the latter two.
243 Although our approach failed to recognize Q9, T15, V20, and Q44, the SVMGA identified the
244 nearby residues, L10, P12, L17, M22, V45, and V49. Residues L10 and P12 also reside at the
245 ice-binding interface.

246 For the detail results obtained for the 369 AFPs in the second independent dataset (Fig. 5), for which
247 no structural information was available, the identification accuracy diminished as the evolutionary
248 distance of a protein sequence increased from the model fish and insect sequences. For sequences
249 with very low SI values (15~20%), especially those from algae, bacteria, and plants, our approach was
250 around 30% accurate. The identification of fish AFPs was around 60% accurate even when sequences
251 with lower than 20% SI values. In fact, we believe that the features encoded in the fish and insect
252 sequences may be used to identify AFPs from evolutionarily divergent organisms. Additionally, as
253 more sequence data for AFPs are accumulated, they can be used to further characterize the
254 mechanisms of cold adaptation. Finally, our approach can be used as an efficient way to obtain high
255 throughput identification of protein function on a genome-wide scale. We have implemented iAFP
256 web service, which is available at <http://140.134.24.89/~iafp/>.

257 **ACKNOWLEDGMENTS**

258 We thank Jenn-Kang Hwang (National Chiao Yung University) for his invaluable comments and
259 crucial insights and Chen-Hsiung Chan (Tzu Chi University) for helpful discussions. This work was
260 supported by grants from the National Science Council, Taiwan to CSY and the National Science
261 Council and China Medical University, Taiwan to CHL. We are grateful for the hardware and software
262 support by the Intelligent Digit Center at Feng Chia University and the Structural Bioinformatics Core
263 Facility at Nation Chiao Tung University, respectively.

264

REFERENCES

1. Fletcher GL, Hew CL, Davies PL (2001) Antifreeze proteins of teleost fishes. *Annu Rev Physiol* 63: 359–390.
2. Knight CA (2000) Structural biology. Adding to the antifreeze agenda. *Nature* 406: 249–251.
3. Fan Y, Liu B, Wang HB, Wang SQ, Wang JF (2002) Cloning of an antifreeze protein gene in carrot and influence on freeze tolerance of transgenic tobaccos. *Plant Cell Rep* 21: 296–301.
4. Rubinsky B, Arav A, Devries AL (1992) The cryoprotective effect of antifreeze glycopeptides from antarctic fishes. *Cryobiology* 29: 69–79.

5. Griffith M, Ewart KV (1995) Antifreeze proteins and their potential use in frozen foods. *Biotechnol Adv* 13: 375–402.
6. Griffith M, Yaish MW (2004) Antifreeze proteins in overwintering plants: a tale of two activities. *Trends Plant Sci* 9: 399–405.
7. Jia Z, Davies PL (2002) Antifreeze proteins: an unusual receptor-ligand interaction. *Trends Biochem Sci* 27: 101–106.
8. Graham LA, Loughheed SC, Ewart KV, Davies PL (2008) Lateral transfer of a lectin-like antifreeze protein gene in fishes. *PLoS ONE* 3: e2616.
9. Doxey AC, Yaish MW, Griffith M, McConkey BJ (2006) Ordered surface carbons distinguish antifreeze proteins and their ice-binding regions. *Nat Biotechnol* 24: 852–855.
10. Graether SP, Sykes BD (2004) Cold survival in freeze-intolerant insects: the structure and function of beta-helical antifreeze proteins. *Eur J Biochem* 271: 3285–3296.
11. Harding MM, Ward LG, Haymet AD (1999) Type I 'antifreeze' proteins. Structure-activity studies and mechanisms of ice growth inhibition. *Eur J Biochem* 264: 653–665.
12. Graether SP, DeLuca CI, Baardsnes J, Hill GA, Davies PL, et al. (1999) Quantitative and qualitative analysis of type III antifreeze protein structure and function. *J Biol Chem* 274: 11842–11847.
13. Graether SP, Kuiper MJ, Gagne SM, Walker VK, Jia Z, et al. (2000) Beta-helix structure and ice-binding properties of a hyperactive antifreeze protein from an insect. *Nature* 406: 325–328.
14. Jia Z, DeLuca CI, Chao H, Davies PL (1996) Structural basis for the binding of a globular antifreeze protein to ice. *Nature* 384: 285–288.
15. Nutt DR, Smith JC (2008) Dual function of the hydration layer around an antifreeze protein revealed by atomistic molecular dynamics simulations. *J Am Chem Soc* 130: 13066–13073.
16. Leinala EK, Davies PL, Jia Z (2002) Crystal structure of beta-helical antifreeze protein points to a general ice binding model. *Structure* 10: 619–627.
17. Fernandez-Recio J, Totrov M, Skorodumov C, Abagyan R (2005) Optimal docking area: a new method for predicting protein-protein interaction sites. *Proteins* 58: 134–143.
18. Lu CH, Chen YC, Yu CS, Hwang JK (2007) Predicting disulfide connectivity patterns. *Proteins* 67: 262–270.
19. Yu CS, Chen YC, Lu CH, Hwang JK (2006) Prediction of protein subcellular localization. *Proteins* 64: 643–651.
20. Yu CS, Lin CJ, Hwang JK (2004) Predicting subcellular localization of proteins for Gram-negative bacteria by support vector machines based on n-peptide compositions. *Protein Sci* 13: 1402–1406.
21. Yu CS, Wang JY, Yang JM, Lyu PC, Lin CJ, et al. (2003) Fine-grained protein fold

- assignment by support vector machines using generalized npeptide coding schemes and jury voting from multiple-parameter sets. *Proteins* 50: 531–536.
22. Wang G, Dunbrack RL, Jr. (2003) PISCES: a protein sequence culling server. *Bioinformatics* 19: 1589–1591.
 23. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, et al. (2000) The Protein Data Bank. *Nucleic Acids Res* 28: 235–242.
 24. Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, et al. (2007) Clustal W and Clustal X version 2.0. *Bioinformatics* 23: 2947–2948.
 25. Bairoch A, Apweiler R (2000) The SWISS-PROT protein sequence database and its supplement TrEMBL in 2000. *Nucleic Acids Res* 28: 45–48.
 26. The Universal Protein Resource (UniProt) in 2010. *Nucleic Acids Res* 38: D142–148.
 27. Hamada T, Ito Y, Abe T, Hayashi F, Guntert P, et al. (2006) Solution structure of the antifreeze-like domain of human sialic acid synthase. *Protein Sci* 15: 1010–1016.
 28. Chen YC, Hwang JK (2005) Prediction of disulfide connectivity from protein sequences. *Proteins* 61: 507–512.
 29. Dubchak I, Muchnik I, Mayor C, Dralyuk I, Kim SH (1999) Recognition of a protein fold in the context of the Structural Classification of Proteins (SCOP) classification. *Proteins* 35: 401–407.
 30. Chang CC, Lin CJ (2001) LIBSVM: a library for support vector machines. pp. Software available from <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
 31. Vapnik V (1995) *The nature of statistical learning theory*. New York Springer.
 32. Matthews BW (1975) Comparison of the predicted and observed secondary structure of T4 phage lysozyme. *Biochim Biophys Acta* 405: 442–451.
 33. Myers EW, Miller W (1988) Optimal alignments in linear space. *Comput Appl Biosci* 4: 11–17.
 34. Rost B (1999) Twilight zone of protein sequence alignments. *Protein Eng* 12: 85–94.
 35. Liou YC, Tocilj A, Davies PL, Jia Z (2000) Mimicry of ice structure by surface hydroxyls and water of a beta-helix antifreeze protein. *Nature* 406: 322–324.
 36. DeLano WL (2002) *The PyMOL Molecular Graphics System* In: Scientific. D, editor. San Carlos, CA, USA. <http://www.pymol.org>.

FIGURE LEGENDS

Fig. 1. Sequence identity distribution for pairs of AFPs. The *x*-axis values are the best pairwise-matched SI values for each AFP sequence against the other 368 sequences. The *y*-axis values are the best pairwise-matched SI values for each of the 369 AFP sequences of the second independent dataset against the 44 sequences of the validation set. A black symbol indicates a correctly identified AFP in the independent data set, and a red symbol indicates an incorrectly identified AFP.

Fig. 2. Examples of key residues mapped onto the surfaces of the eight representative AFPs used in the cross-validation tests. The structures were drawn with PyMOL [36]. The residues colored in gray were not identified as key residues. The residues in red obtained more votes than did the residues in yellow. (a) 1c3y; (b) 1ezg; (c) 1eww; (d) 2pne; (e) 1c89; (f) 1wfa; (g) 2py2; (h) 2afp.

Fig. 3. Difference of the number of SVMGA key residues extracted from the 44 AFP and 3762 non-AFP sequences in cross-validation dataset, respectively. Each black bar represents the mean \pm standard deviations of coverage percentage a SVMGA residue was included in a AFP sequence. Each gray bar represents for non-AFP sequence.

Fig. 4. The surface of the eelpout type III AFP (PDB ID 1msi) drawn with PyMOL [36]. (a) The key residues selected by the SVMGA are labeled in black words. Residues Q9 and N14, which were identified as key residues in a mutagenesis study but not by the SVMGA, are labeled in blue. (b) A view of the ice-binding interface, wherein all residues that are part of the interface are labeled. The residues identified by SVMGA are in red and yellow. Residues known to be important in ice binding, but not identified by the SVMGA, are in cyan. Residue I13, which was not identified by the SVMGA, is in gray. Its status as a key residue has not been determined by a mutagenesis study.

Fig. 5. The identification accuracy for the 369 AFPs from the second independent set. Each bar correlates the identification accuracy with a range of maximum SI values, which was found using the *y* axis of Figure 1 in detail ranges of SI for different species.

Table 1. The eight protein subsets used for cross-validation testing.

Subset	Type	PDB ID
1	insect AFP	1c3y
2	Type III fish AFP	1c89 ; 3nla; 1ucs; 1ops; 1kde; 1ame; 1msi; 1b7i; 1b7j; 1b7k; 1ekl; 1gzi; 1hg7; 1jab; 1msj; 2ame; 2jia; 2msi; 2msj; 2spg; 3ame; 3msi; 4ame; 4msi; 5msi; 6ame; 6msi; 7ame; 7msi; 8ame; 8msi; 9ame; 9msi;
3	β -helical insect AFP	1ezg
4	Type I fish AFP	1wfa ; 1j5b; 1y03
5	β -helical insect AFP	1eww ; 110s; 1m8n
6	insect AFP	2pne
7	Type II fish AFP	2py2
8	Type II fish AFP	2afp

Notes: The sequences of the PDB codes given in bold type were used for the single representative AFP mode.

Table 2. The number of antifreeze protein sequences for a given type of organism in the independent dataset that contained 369 AFPs.

Organism	Number of sequences
Algae	17
Bacteria	101
Fish	123
Insects	105
Plants	23

Table 3. The performances of SVM and SVMGA for the eight-fold cross-validation tests that used the single representative AFP mode or the multiple representative AFP mode.

Number	Subset	SVM		SVMGA		§§Doxey et al.[9]
		C+X ₃ +V ₃ X ₅		§14 Feature Schemes		
1 (1)	1	0	(0)	1	(1)	-
1 (33)	2	0	(0)	1	(15)	(3)
1 (1)	3	1	(1)	1	(1)	(1)
1 (3)	4	0	(0)	1	(1)	(3)
1 (3)	5	1	(2)	1	(3)	(2)
1 (1)	6	1	(1)	1	(1)	-
1 (1)	7	1	(1)	1	(1)	-
1 (1)	8	1	(1)	1	(1)	(0)
AFP accuracy		62.5%	(13.6%)	100.0%	(54.5%)	(90.0%)
AFP precision		21.7%	(25.0%)	10.4%	(25.8%)	(42.9%)
Overall accuracy		99.4%	(98.5%)	98.2%	(97.7%)	(99.6%)
MCC		0.367	(0.178)	0.319	(0.365)	(0.620)
TP		5	(6)	8	(24)	(9)
TN		3744	(3744)	3693	(3693)	(3184)
FP		18	(18)	69	(69)	(12)
FN		3	(38)	0	(20)	(1)

Notes: Values given in parentheses are the number of homologous proteins accurately recognized using in the multiple representative AFP mode.

§14 feature schemes: $\sum_{k=1}^3 X_k^{A1} + \sum_g D_g + \sum_S X_{k=5}^S + \sum_{l \in S'} W_l$ where $g = \{0,1,2,3,5\}$, $S = \{H_3, V_3, P_3, S_2\}$, and $S' = \{9,15\}$

§§Doxey and colleagues used structure as the property to identify 10 AFPs in their dataset excellently. Only 2atp, for which its NMR structure was used, was not identified correctly.

Table 4. The feature schemes that enabled the recognition of the AFP in a subset when the single representative mode was used. The filled circles correlate the feature schemes with the AFPs that they identified. The AFPs are denoted according to their subsets.

Subset	Feature Scheme									
	C	W_1	D_0	D_2	D_3	S_2X_5	H_3X_5	P_3X_5	V_3X_5	Z_3X_5
1			•							
2		•	•			•				
3	•	•	•	•	•	•			•	•
4			•				•	•	•	
5	•	•		•	•	•	•	•	•	
6							•	•	•	•
7	•	•	•		•		•			
8	•	•	•			•	•	•		

Table 5. An example of votes acquired by residues in a sequence from 1msi.

Sequence	Q ⁹	L ¹⁰	I ¹¹	P ¹²	I ¹³	N ¹⁴	T ¹⁵	A ¹⁶	L ¹⁷	T ¹⁸
Coding	<i>C</i>	*	*		*					*		
	<i>X₂</i>				*			*	*		*	
	<i>X₃</i>			*		*						
	<i>D₀</i>									*	*	
	<i>D₁</i>		*									
	<i>D₂</i>							*			*	
	<i>D₃</i>											
	<i>D₅</i>										*	
	<i>O₃X₅</i>								*	*	*	
	<i>P₃X₅</i>							*	*	*		
	<i>V₃X₅</i>		*	*	*				*	*	*	
	<i>S₂X₅</i>											
	<i>W₉</i>		*	*	*	*	*		*	*		
	<i>W₁₅</i>											
Votes	1	4	3	4	2	1	3	5	6	6

Figure 1
[Click here to download high resolution image](#)

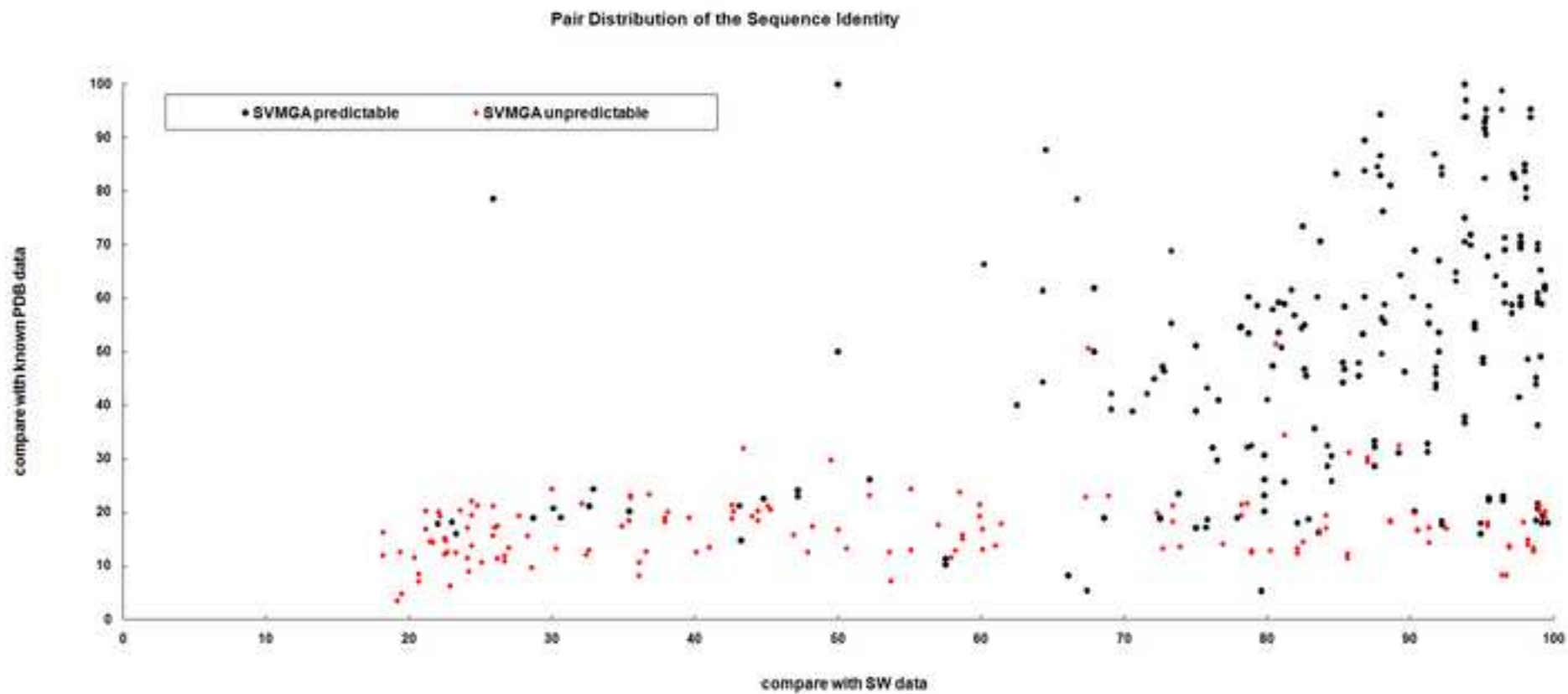


Figure 3
[Click here to download high resolution image](#)

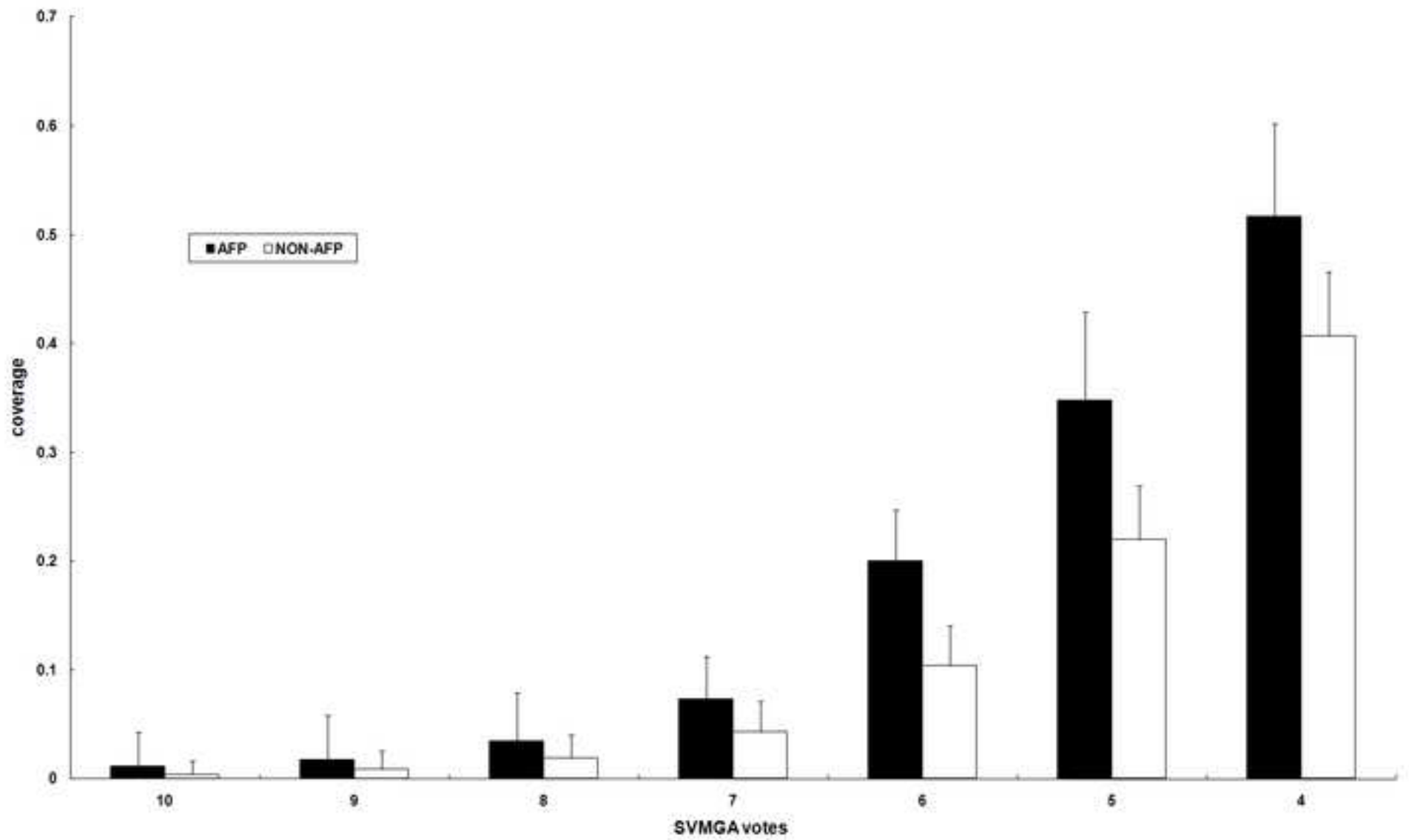


Figure 4
[Click here to download high resolution image](#)

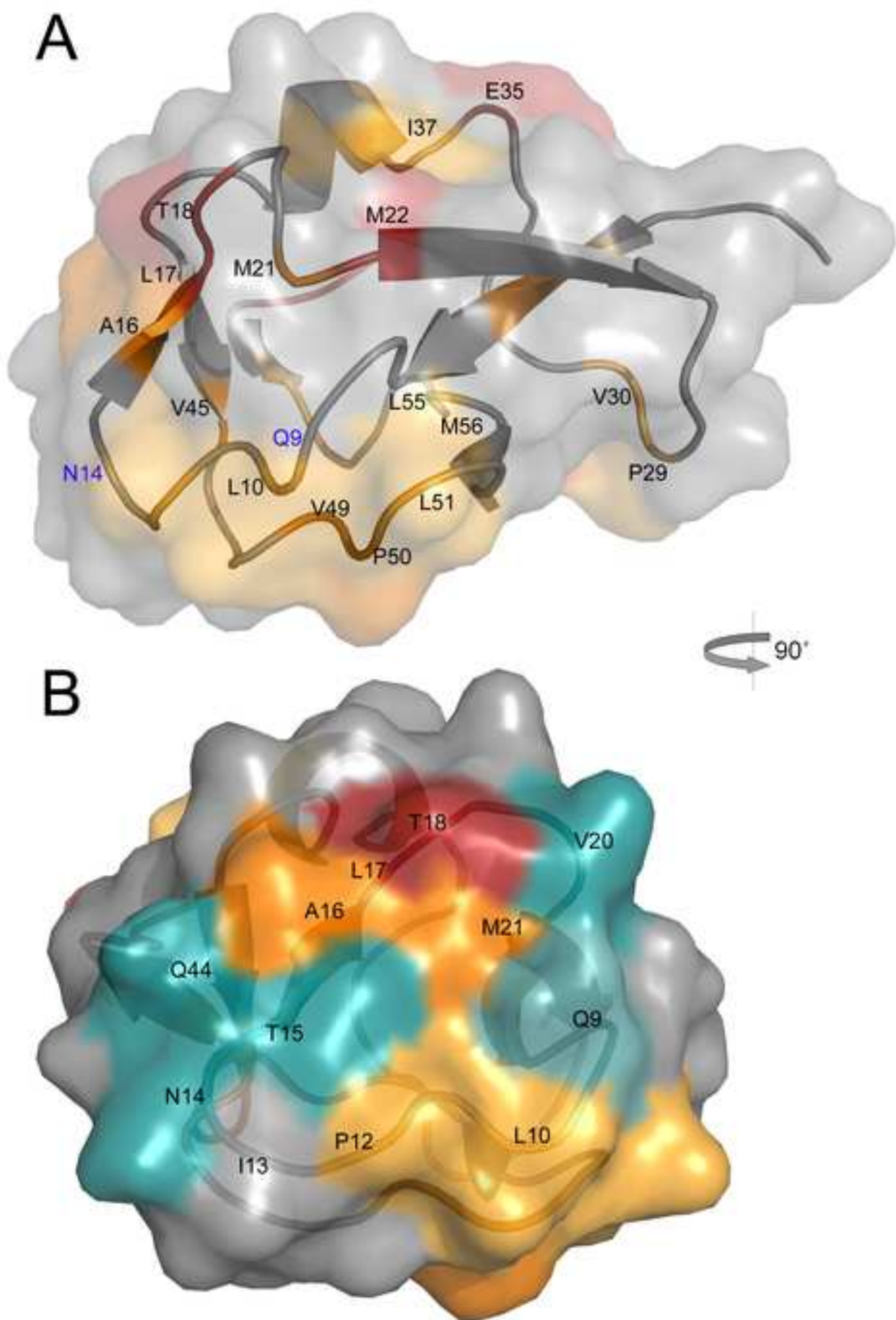
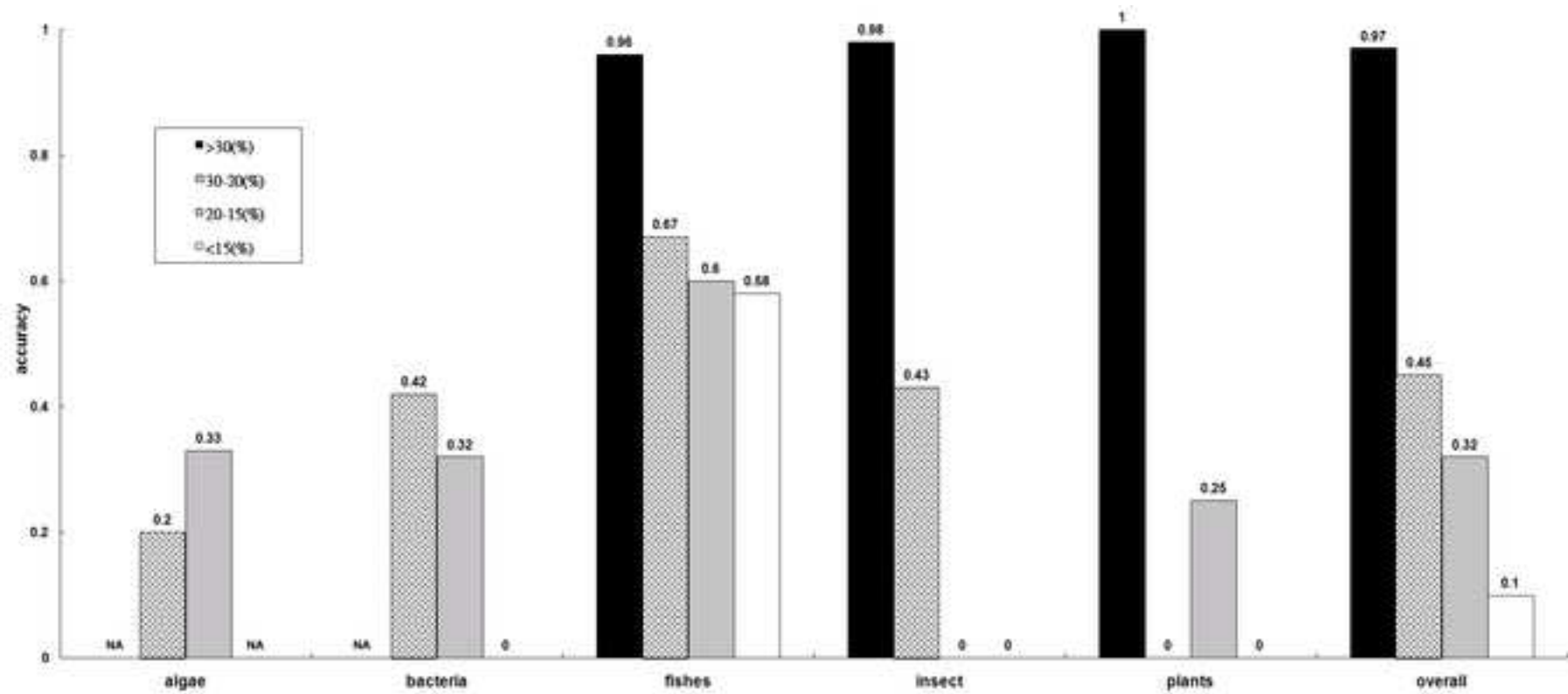


Figure 5
[Click here to download high resolution image](#)



independent dataset - 369 AFP list

[Click here to download Supporting Information: AFP369SW.pdf](#)